

Vapnik-Chervonenkis Dimension of Axis-Parallel Cuts

Servane Gey*

February 24, 2013

Abstract

The Vapnik-Chervonenkis (VC) dimension of the set of half-spaces of \mathbb{R}^d with frontiers parallel to the axes is computed exactly. It is shown that it is much smaller than the intuitive value of d . A good approximation based on the Stirling's formula proves that it is more likely of the order $\log_2 d$.

This result may be used to evaluate the performance of classifiers or regressors based on dyadic partitioning of \mathbb{R}^d for instance. Algorithms using axis-parallel cuts to partition \mathbb{R}^d are often used to reduce the computational time of such estimators when d is large.

Keywords: Vapnik-Chervonenkis dimension, axis-parallel cuts.

MSC 2010 Classification: 62G99 62H99

1 Introduction

The VC dimension of a set of subsets has been introduced by Vapnik and Chervonenkis [9, 10] to measure its complexity. The VC dimension of a real-valued function space \mathcal{F} is then the VC dimension of $\{\{x; f(x) \geq 0\}; f \in \mathcal{F}\}$. In particular, the VC dimension of sets of classifiers or regressors appears commonly in the statistical learning area when evaluating their performance. For example, Vapnik's theory in the classification framework is now widely known (see [3] for instance): let (X, Y) be a couple of variables taking values in $\mathbb{R}^d \times \{0; 1\}$, and let \mathcal{L} be a sample of n independent replications of (X, Y) . If \hat{f} is a classifier minimizing the average misclassification rate of \mathcal{L} on a set

*Laboratoire MAP5 - UMR 8145, Université Paris Descartes, 75270 Paris Cedex 06, France - Servane.Gey@parisdescartes.fr

of classifiers having finite VC dimension V , then, without further assumption on the distribution P of (X, Y) , the performance of \hat{f} is evaluated as follows:

$$\mathbb{E}_{\mathcal{L}} \left[P \left(\hat{f}(X) \neq Y \right) \right] \leq C_1 \text{bias}^2(\hat{f}) + C_2 \sqrt{\frac{V}{n}}, \quad (1)$$

where $\mathbb{E}_{\mathcal{L}}$ denotes the expectation with respect to the sample distribution, $\text{bias}(\hat{f})$ denotes the bias of the classifier \hat{f} , and C_1 and C_2 are absolute constants.

Functional estimates defined on partitions of \mathbb{R}^d are often used to estimate relationships between two variables $X \in \mathbb{R}^d$ and $Y \in \{0; 1\}$ or $Y \in \mathbb{R}$ (such as histograms, piecewise polynomials, or splines for example). In many cases, the VC dimension of the set of subsets used to construct the partition appears inside risk bounds when evaluating the performance of such estimators. For example, if the set used is the set of all half-spaces of \mathbb{R}^d , often its VC dimension $d + 1$ has to be taken into account.

When d is large, it is often computationally easier to construct partitions using axis-parallel cuts. For example, some theoretical developments on dyadic partitions of \mathbb{R}^2 are given in [4, 1], and the VC dimension of axis-parallel cuts appears more particularly in the results obtained on the performance of classification and regression binary decision trees (CART) introduced by Breiman *et. al* [2] in 1984, and theoretically studied in [8, 7, 5, 6].

2 Reminder about VC Dimension

The VC dimension of a set \mathcal{A} of subsets of some measurable space \mathcal{X} is based on counting the number of intersects of \mathcal{A} with a finite set of fixed points in \mathcal{X} .

Definition 1 (Vapnik-Chervonenkis Dimension). *Let \mathcal{A} be a set of subsets of some measurable space \mathcal{X} . Then $(x_1, \dots, x_n) \in \mathcal{X}^n$ will be said to be shattered by \mathcal{A} if all subsets of $\{x_1; \dots; x_n\}$ are covered by \mathcal{A} , that is if $|\{\{x_1, \dots, x_n\} \cap A ; A \in \mathcal{A}\}| = 2^n$.*

The Vapnik-Chervonenkis dimension $VC(\mathcal{A})$ of \mathcal{A} is then defined as the maximal integer n such that there exists n points in \mathcal{X} shattered by \mathcal{A} , i.e.

$$VC(\mathcal{A}) = \max \left\{ n ; \max_{(x_1, \dots, x_n) \in \mathcal{X}^n} |\{\{x_1, \dots, x_n\} \cap A ; A \in \mathcal{A}\}| = 2^n \right\}.$$

If no such n exists, then $VC(\mathcal{A}) = +\infty$.

Thus, it is easily seen that the larger $VC(\mathcal{A})$, the more complex \mathcal{A} .

For example, if $\mathcal{A} = \{] - \infty; x] ; x \in \mathbb{R} \}$, then $VC(\mathcal{A}) = 1$; or if \mathcal{A} is the set of all half-spaces in \mathbb{R}^d , then $VC(\mathcal{A}) = d + 1$.

Since axis-parallel cuts is a subset of the set of all half-spaces in \mathbb{R}^d , it could be natural to think that its VC dimension is of order d . Actually, it is shown in what follows that it is of order $\log_2 d$.

3 VC Dimension of axis-parallel cuts

We give a formula to compute the VC dimension of axis-parallel cuts in \mathbb{R}^d . Since the obtained formula is not always easy to handle, an approximation is also given.

Lemma 1. *Let*

$$\mathcal{A}_d = \left\{ \{x \in \mathbb{R}^d ; x^i \leq a\} ; i = 1, \dots, d, a \in \mathbb{R} \right\}.$$

Then

$$VC(\mathcal{A}_d) = \max \left\{ n ; \binom{n}{\lfloor n/2 \rfloor} \leq d \right\},$$

where $\lfloor x \rfloor$ denotes the integer part of x .

Furthermore, the following approximation of $VC(\mathcal{A}_d)$ is available for all $d \geq 2$:

$$\frac{\log d}{\log 2} - 0.38 \leq VC(\mathcal{A}_d) \leq \frac{\log(d\sqrt{d+3})}{\log 2} + 0.51.$$

Figure 1 shows that $VC(\mathcal{A}_d)$ is a piecewise constant function of the space dimension d , which increases at a rate much smaller than the intuitive value of d . It also shows that the bounds computed from the Stirling's formula are sharp.

Proof. The idea is that, to have n points (x_1, \dots, x_n) shattered by \mathcal{A}_d , all the subsets of $\{x_1, \dots, x_n\}$ should be covered by \mathcal{A}_d . But, if there exists $p \leq n$ such that there is more than $d + 1$ subsets of $\{x_1, \dots, x_n\}$ having p elements, then \mathcal{A}_d will miss at least $\binom{n}{p} - d$ subsets: let $n \geq 1$ and (x_1, \dots, x_n) be n points in \mathbb{R}^d . Suppose that n is such that $\binom{n}{\lfloor n/2 \rfloor} > d$. This means that there are at least $d + 1$ subsets of $\{x_1, \dots, x_n\}$ of size $\lfloor n/2 \rfloor$. For each

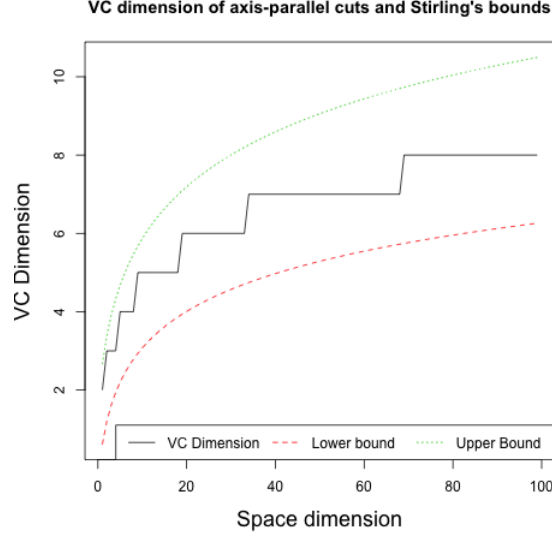


Figure 1: $VC(\mathcal{A}_d)$ with respect to the space dimension d and Stirling's bounds.

coordinate $i = 1, \dots, d$, let us denote by $x_{i(\cdot)}$ the ordered statistic computed from the i^{th} coordinate of (x_1, \dots, x_n) , that is, for all $i = 1, \dots, d$,

$$x_{i(1)}^i \leq x_{i(2)}^i \leq \dots \leq x_{i(n)}^i.$$

Let $p = \lfloor n/2 \rfloor$ and let

$$\begin{aligned} \mathcal{B}_p &= \{ \{x_{i(1)}; \dots; x_{i(p)}\} ; i = 1, \dots, d \text{ and } |\{x_{i(1)}; \dots; x_{i(p)}\}| = p \}, \\ \mathcal{B}_p^c &= \{ B \subset \{x_1, \dots, x_n\}; |B| = p \text{ and } B \notin \mathcal{B}_p \}. \end{aligned}$$

Hence \mathcal{B}_p is covered by \mathcal{A}_d (by simply taking $A = \{x^i \leq (x_{i(p)}^i + x_{i(p+1)}^i)/2\}$ for each coordinate), and we have that:

$$|\mathcal{B}_p| \leq d \text{ and } |\mathcal{B}_p^c| \geq \binom{n}{p} - d > 0.$$

Let $B \in \mathcal{B}_p^c$ and $A = \{x^i \leq a\} \in \mathcal{A}_d$. If $|\{x_1, \dots, x_n\} \cap A| \neq p$, then $\{x_1, \dots, x_n\} \cap A \neq B$. Else, since $\{x_1, \dots, x_n\} \cap A = \{x_j ; x_j^i \leq a\}$, we have that $x_{i(j)}^i \leq a$ for all $j = 1, \dots, p$, and $x_{i(j)}^i > a$ for all $j = p+1, \dots, n$. So $\{x_1, \dots, x_n\} \cap A = \{x_{i(1)}; \dots; x_{i(p)}\}$ and $|\{x_{i(1)}; \dots; x_{i(p)}\}| = p$, leading to $\{x_1, \dots, x_n\} \cap A \in \mathcal{B}_p$, and then to $\{x_1, \dots, x_n\} \cap A \neq B$. So, for all $B \in \mathcal{B}_p^c$

and all $A \in \mathcal{A}_d$, $\{x_1, \dots, x_n\} \cap A \neq B$.

So, if $\binom{n}{\lfloor n/2 \rfloor} > d$, (x_1, \dots, x_n) can not be shattered by \mathcal{A}_d . Thus

$$VC(\mathcal{A}_d) \leq \max \left\{ n ; \binom{n}{\lfloor n/2 \rfloor} \leq d \right\}.$$

Let $n \geq 1$ such that $\binom{n}{\lfloor n/2 \rfloor} \leq d$. Let (x_1, \dots, x_n) be n points of \mathbb{R}^d defined as follows: for each coordinate $i = 1, \dots, \binom{n}{\lfloor n/2 \rfloor}$, let $\{i_1; \dots; i_{\lfloor n/2 \rfloor}\}$ be the i^{th} subset of $\lfloor n/2 \rfloor$ indices in $\{1; \dots; n\}$, where the indices are denoted in ascending order, i.e.:

$$1 \leq i_1 < \dots < i_{\lfloor n/2 \rfloor} \leq n.$$

Since $\binom{n}{\lfloor n/2 \rfloor} \leq d$, we obtain $\binom{n}{\lfloor n/2 \rfloor}$ distinct subsets of indices.

Hence we take for each such coordinate

$$x_{i_k}^i = k.$$

Then the remaining values of (x_1, \dots, x_n) are taken as follows:

- Since $\binom{n}{\lfloor n/2 \rfloor + 1} \leq d$, for each subset $\{i_1; \dots; i_{\lfloor n/2 \rfloor + 1}\}$ of $\{1; \dots; n\}$ with $\lfloor n/2 \rfloor + 1$ elements, there exists $i' \in \{1; \dots; \binom{n}{\lfloor n/2 \rfloor}\}$ such that $\{i_1; \dots; i_{\lfloor n/2 \rfloor}\} = \{i'_1; \dots; i'_{\lfloor n/2 \rfloor}\}$. Then take $x_{i_{\lfloor n/2 \rfloor + 1}}^{i'} = \lfloor n/2 \rfloor + 1$. Let us note that, if n is odd, there is a bijection between i and i' .
- Let $\{j_1; \dots; j_m\} = \{j \notin \{i_1; \dots; i_{\lfloor n/2 \rfloor + 1}\}\}$, with $j_1 < \dots < j_m$, and let $j_0 = i_{\lfloor n/2 \rfloor + 1}$. Then take $x_{j_k}^{i'} = x_{j_{k-1}}^{i'} + 1$.

If not filled, the last coordinates are set to be equal to n .

Hence, we obtain that, for all $j \notin \{i_1; \dots; i_{\lfloor n/2 \rfloor}\}$, $x_j^i \geq \lfloor n/2 \rfloor + 1$.

Then (x_1, \dots, x_n) is shattered by \mathcal{A}_d : for $p \in \{0; \dots; n\}$, let $B = \{x_{i_1}; \dots; x_{i_p}\} \subset \{x_1, \dots, x_n\}$, with $1 \leq i_1 < i_2 < \dots < i_p \leq n$ as soon as $p \neq 0$.

If $p = 0$, let

$$i_0 = \operatorname{argmin}_{1 \leq i \leq d} \min_j x_j^i,$$

and take $A = \{x^{i_0} \leq \min_j x_j^{i_0} - 1\}$. Then $B = \{x_1, \dots, x_n\} \cap A = \emptyset$.

If $p = n$, let

$$i_n = \operatorname{argmax}_{1 \leq i \leq d} \max_j x_j^i,$$

and take $A = \{x^{i_n} \leq \max_j x_j^{i_n} + 1\}$. Then $B = \{x_1, \dots, x_n\} \cap A = \{x_1, \dots, x_n\}$.

If $0 < p \leq \lfloor n/2 \rfloor$, let $A \in \mathcal{A}_d$ be the subset defined by $A = \{x^i \leq p + 1/2\}$, with i the coordinate corresponding to a subset of indices $\{i_1; \dots; i_{\lfloor n/2 \rfloor}\}$ containing $\{i_1; \dots; i_p\}$. Then, by definition of (x_1^i, \dots, x_n^i) , $B = \{x_1, \dots, x_n\} \cap A$.

If $\lfloor n/2 \rfloor + 1 \leq p < n$, let i' be the coordinate corresponding to the configuration $\{i_1; \dots; i_{\lfloor n/2 \rfloor + 1}\}$ (as defined by (x_1, \dots, x_n)). Let $A \in \mathcal{A}_d$ be the subset defined by $A = \{x^{i'} \leq p + 1/2\}$. Then, by definition of $(x_1^{i'}, \dots, x_n^{i'})$, $B = \{x_1, \dots, x_n\} \cap A$.

Thus

$$VC(\mathcal{A}_d) \geq \max \left\{ n ; \binom{n}{\lfloor n/2 \rfloor} \leq d \right\}.$$

Then, the lower and upper bounds of $VC(\mathcal{A}_d)$ are computed by using the Stirling's formula: for all $n \geq 1$ we have

$$\sqrt{2\pi}e^{-(n+1)}(n+1)^{n+\frac{1}{2}} \leq n! \leq \sqrt{2\pi}e^{-(n+1)}e^{\frac{1}{12(n+1)}}(n+1)^{n+\frac{1}{2}}.$$

A simple calculation gives the following: if n is even, then

$$\binom{n}{n/2} \leq e^{\frac{1}{12(n+1)}} \frac{e}{\sqrt{2\pi}} 2^{n+1} \frac{(n+1)^{n+1/2}}{(n+2)^{n+1}} \leq \frac{e^{1+\frac{1}{36}}}{\sqrt{6\pi}} 2^{n+1},$$

and if n is odd, then

$$\binom{n}{\lfloor n/2 \rfloor} \leq e^{\frac{1}{12(n+1)}} \frac{e}{\sqrt{2\pi}} 2^{n+1} \frac{(n+1)^{(n+1)/2}}{(n+3)^{n/2+1}} \leq \frac{e^{1+\frac{1}{24}}}{\sqrt{2\pi}} 2^n.$$

Thus, if

$$\frac{e^{1+\frac{1}{24}}}{\sqrt{6\pi}} 2^{n+1} \leq d,$$

then $\binom{n}{\lfloor n/2 \rfloor} \leq d$. Taking the logarithm leads to the lower bound of $VC(\mathcal{A}_d)$.

On the other hand, if $\binom{n}{\lfloor n/2 \rfloor} \leq d$, we have that, if n is even,

$$\binom{n}{n/2} \geq e^{-\frac{1}{3n+6}} \frac{e}{\sqrt{2\pi}} 2^{n+1} \frac{(n+1)^{n+1/2}}{(n+2)^{n+1}} \geq \frac{e^{-\frac{1}{12}}}{\sqrt{2\pi}} 2^{n+1} \frac{1}{\sqrt{d+2}},$$

and if n is odd,

$$\binom{n}{\lfloor n/2 \rfloor} \geq e^{-\frac{n+2}{3(n+3)(n+1)}} \frac{e}{\sqrt{2\pi}} 2^{n+1} \frac{(n+1)^{(n+1)/2}}{(n+3)^{n/2+1}} \geq \frac{e^{-\frac{1}{8}}}{\sqrt{2\pi}} 2^{n+1} \frac{1}{\sqrt{d+3}}.$$

Thus, since, for all n such that $\binom{n}{\lfloor n/2 \rfloor} \leq d$, $\frac{e^{1-\frac{1}{8}}}{\sqrt{2\pi}} 2^{n+1} \leq d$, the upper bound of $VC(\mathcal{A}_d)$ is found by taking the logarithm of this last expression. \square

References

- [1] AKAKPO, N. Adaptation to anisotropy and inhomogeneity via dyadic piecewise polynomial selection. *Mathematical Methods of Statistics* 21, 1 (2012), 1–28.
- [2] BREIMAN, L., FRIEDMAN, J. H., OLSHEN, R. A., AND STONE, C. J. *Classification And Regression Trees*. Chapman & Hall, 1984.
- [3] DEVROYE, L., GYÖRFI, L., AND LUGOSI, G. *A probabilistic theory of pattern recognition*, vol. 31 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 1996.
- [4] DONOHO, D. L. CART and best-ortho-basis : A connection. *The Annals of Statistics* 25, 5 (1997), 1870–1911.
- [5] GEY, S. Risk bounds for cart classifiers under a margin condition. *Pattern Recognition* 45 (2012), 3523–3534.
- [6] GEY, S., AND MARY HUARD, T. Risk bounds for embedded variable selection in classification trees. Tech. rep., arxiv, 1108.0757v1, 2011.
- [7] GEY, S., AND NEDELEC, E. Model selection for CART regression trees. *IEEE Trans. Inform. Theory* 51, 2 (2005), 658–670.
- [8] NOBEL, A. B. Analysis of a complexity-based pruning scheme for classification trees. *IEEE Trans. Inform. Theory* 48, 8 (2002), 2362–2368.
- [9] VAPNIK, V. N., AND CHERVONENKIS, A. Y. Theory of uniform convergence of frequencies of events to their probabilities and problems of search for an optimal solution from empirical data. *Avtomat. i Telemekh.*, 2 (1971), 42–53.

- [10] VAPNIK, V. N., AND CHERVONENKIS, A. Y. *Teoriya raspoznavaniya obrazov. Statisticheskie problemy obucheniya*. Izdat. “Nauka”, Moscow, 1974.